

1 Population size history from short genomic scaffolds: how short 2 is too short?

3
4 Graham Gower^{*1}, Jono Tuke^{2,3}, AB Rohrlach^{2,3}, Julien Soubrier^{1,4}, Bastien Llamas¹, Nigel Bean^{2,3},
5 and Alan Cooper¹

6 ¹Australian Centre for Ancient DNA, School of Biological Sciences, The Environment Institute, The
7 University of Adelaide, Adelaide, South Australia 5005, Australia

8 ²School of Mathematical Sciences, The University of Adelaide, Adelaide, South Australia 5005,
9 Australia

10 ³ARC Centre of Excellence for Mathematical and Statistical Frontiers, The University of Adelaide,
11 Adelaide, South Australia 5005, Australia

12 ⁴Genetics and Molecular Pathology, SA Pathology, Adelaide, South Australia 5000, Australia

13 ^{*}Corresponding author: graham.gower@adelaide.edu.au

14 Abstract

15 The Pairwise Sequentially Markov Coalescent (PSMC), and its extension PSMC', model past
16 population sizes from a single diploid genome. Both models have been widely applied, even to
17 organisms with scaffold-level genome reference assemblies of limited contiguity. However it is unclear
18 how PSMC and PSMC' perform on short scaffolds. We evaluated `psmc` and `msmc`, implementations of
19 the PSMC and PSMC' models respectively, on simulated genomes with low contiguity, and compared
20 results to those from fully contiguous data. Simulations with scaffolds from 100 Mb to 10 kb revealed
21 that `psmc` maintains high accuracy down to lengths of 100 kb, while `msmc` is accurate down to 1 Mb.
22 The discrepancy is not due to differing models, but stems from an implementation detail of
23 `msmc`—homozygous tracts at the ends of scaffolds are discarded, making `msmc` unreliable for low
24 contiguity genomes. We recommend excluding data that are aligned to shorter scaffolds when
25 undertaking demographic inference.

26 Introduction

27 The process of joining (coalescing) and splitting (recombining) lineages backwards-in-time for a sample of
28 homologous sequences is described by the coalescent with recombination (Hudson, 1990). An important
29 consequence of recombination is that there can be many distinct genealogies, known as marginal genealo-
30 gies, at different locations along the sequence (Griffiths and Marjoram, 1997). The sequentially Markov
31 coalescent (SMC, McVean and Cardin (2005)) models recombination as a Poisson process left-to-right
32 along the sequence, approximating the coalescent with recombination by treating the marginal genealogy
33 on the right of a recombination as a modification of the marginal genealogy on the left of the recom-
34 bination. In this sense the approximation is a Markovian process along the sequence, and substantially
35 reduces model complexity for long sequences compared to the full coalescent with recombination (Wiuf
36 and Hein, 1999).

37 The Pairwise Sequentially Markov Coalescent (PSMC) uses a special case of the SMC approxima-
38 tion, restricted to pairs of sequences, to estimate the distribution of coalescent times within a single
39 diploid genome (Li and Durbin, 2011). PSMC scans along a contiguous segment of the genome and
40 considers marginal genealogies, using their distinct pairwise coalescent times as the unknown states in a
41 hidden Markov model (HMM). To enable parameter estimation, continuous time is approximated by a
42 finite partition of time intervals, and transition probabilities are inferred by Baum-Welch iteration of the
43 forward-backward algorithm. Each genotype at consecutive genomic coordinates provides a new observa-
44 tion for the HMM, a homozygote or a heterozygote, with their emission probabilities determined by the
45 pairwise coalescent time at the current locus, and the genome-wide mutation rate. The population size
46 in a given time interval is inversely proportional to the rate of coalescence, as inferred by maximising the
47 fit of the model to both the HMM transition matrix and the emission probabilities.

48 The Multiple Sequential Markov Coalescent (MSMC, Schiffels and Durbin (2014)) is an extension to
49 PSMC, and models the distribution of first-coalescent times of two or more haploid sequences. If used
50 with only two haploid sequences, MSMC closely matches the PSMC model, with the exception that it
51 implements SMC' (Marjoram and Wall, 2006), a refinement of SMC incorporating recombinations that
52 immediately coalesce back to the same lineage. For this reason the MSMC model, when applied to a
53 diploid genome, is referred to as PSMC'. Compared to PSMC, the genome wide recombination rate
54 is more accurately estimated under the PSMC' model, but population size estimates are qualitatively
55 similar (Schiffels and Durbin, 2014).

56 Other approaches for inferring population size histories typically require either phased genotypes,
57 multiple individuals, or both (Dutheil *et al.*, 2009; Gutenkunst *et al.*, 2009; Sheehan *et al.*, 2013; Boitard
58 *et al.*, 2016; Terhorst *et al.*, 2017). However, in small scale studies of non-model organisms, it is common
59 for only one individual, or a few individuals, from a single population to be sequenced, and genotypes are
60 unlikely to be phased. Population size history, particularly in the recent past, can also be estimated from
61 the length distribution of tracts of identity-by-descent (Palamara *et al.*, 2012), identity-by-state (Harris
62 and Nielsen, 2013), or runs of homozygosity (MacLeod *et al.*, 2013). While potentially useful for a single
63 diploid individual, such approaches are not readily applicable to short scaffolds, where such tracts may

64 be broken across scaffold boundaries. In contrast, PSMC and PSMC' are very attractive as they require
65 only diploid genotypes for a single individual, which need not be phased.

66 By using the sequentially Markovian approximation, PSMC and derived methods implicitly assume
67 that genomic information is contiguous. While initially applied to human datasets, which have very
68 high contiguity, PSMC and PSMC' have since been applied to many non-model organisms where the
69 contiguity of genomic sequences may be poor (Zhao *et al.*, 2013; Dobrynin *et al.*, 2015; Mays *et al.*, 2018;
70 Kozma *et al.*, 2016; Feigin *et al.*, 2018). In particular, demographic history is regularly inferred from a
71 *de novo* assembly as part of genome sequencing projects. Due to time and funding constraints, genome
72 assemblies are often constructed from only short read sequencing data, and assembled into contigs or
73 short scaffolds. These cannot be ordered or oriented with respect to one another (violating the SMC
74 model), nor anchored to physical chromosomes. Where sequencing data is aligned to such assemblies, the
75 genomic information used for population size inference inherits the low contiguity of the assembly. While
76 small gaps in coverage along a scaffold can be handled gracefully, the HMM must be applied separately
77 to each distinct scaffold, and it is not clear what the length threshold is to obtain robust population size
78 inferences.

79 Results and Discussion

80 Simulations

81 To assess the impact of reference genome contiguity on population size estimates, we simulated genomes
82 for populations with three different demographic histories: a constant population size; a bottleneck; and
83 recovery following a bottleneck (Fig. 1A). For each demographic scenario, we simulated 10 independent
84 populations and sampled 20×100 Mb haploid chromosomes, representing 10 diploid genomes from each
85 population. New datasets were then created by fragmenting each genome into equally sized scaffolds at
86 four distinct lengths, 10 Mb, 1 Mb, 100 kb, and 10 kb. Population size histories were then inferred for all
87 fragmented and unfragmented datasets using `psmc` (Li and Durbin, 2011) and `msmc` (Schiffels and Durbin,
88 2014), implementations of PSMC and PSMC' respectively.

89 Mean squared error

90 In measuring the error of estimates, Li and Durbin (2011) compared population size inferences to the
91 values that were simulated, but excluded time intervals in the recent and distant past. Population
92 size estimates are expected to be unreliable for times outside a certain range since a typical genome
93 contains relatively few breakpoints corresponding to recombination events in the very recent or very
94 distant past. However, excluding temporal intervals requires advance knowledge of where the method
95 may lose resolution, and this is dependent upon the population size history itself.

96 To quantify estimation error, we used inferences from the unfragmented datasets as the 'truth', not
97 the values that were simulated. A loess smooth function (Cleveland *et al.*, 1992) was fitted to the
98 unfragmented inferences for each simulated population, separately for `psmc` and `msmc`, using population

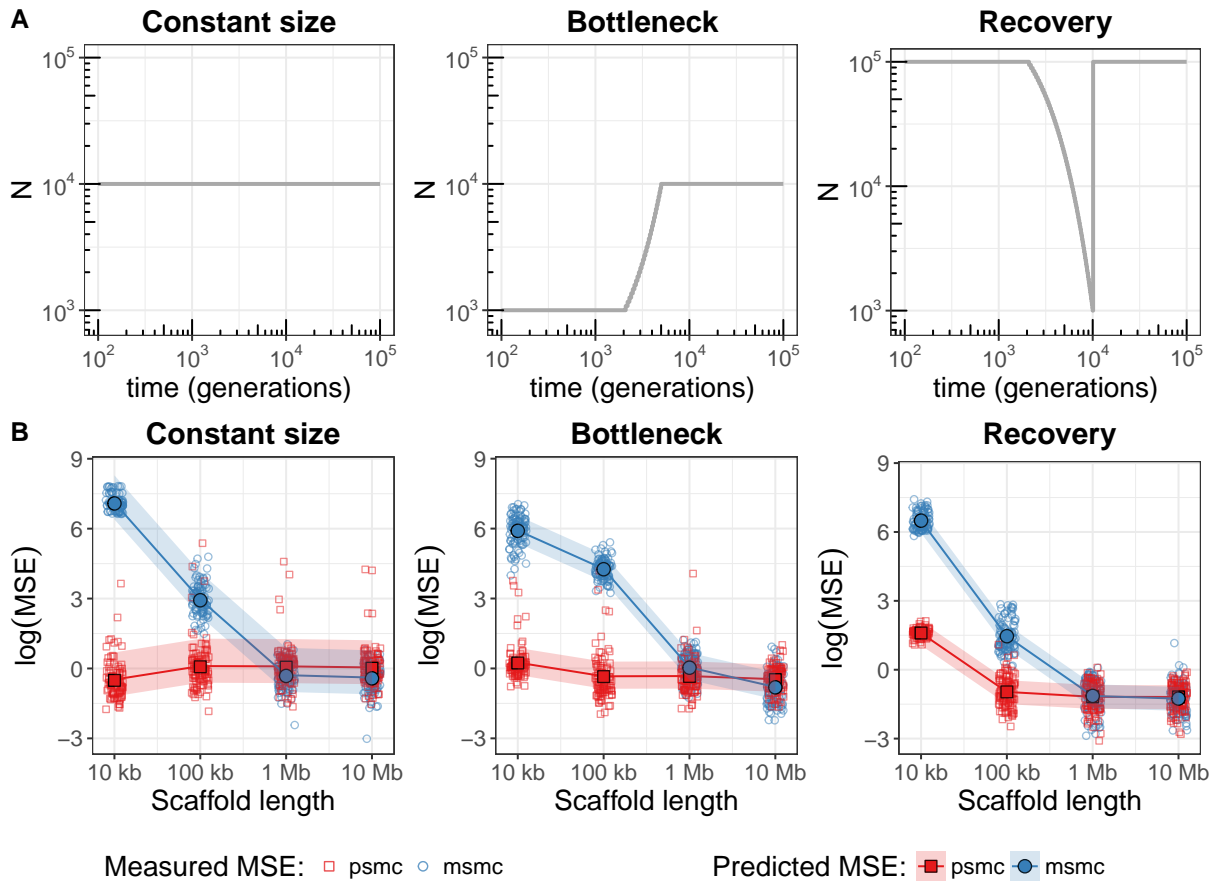


Figure 1: **A)** Simulated population size histories. **B)** Mean squared error (MSE) of population size inferences from simulations shown immediately above. Larger values indicate a loss of accuracy in the population size estimate. Small hollow markers indicate MSE for distinct simulated individuals (100 Mb per individual; 10 individuals each from 10 populations), with red squares for psmc and blue circles for msmc. Data from each simulated individual was artificially fragmented to emulate genome sequences aligned to a scaffold-level reference assembly. At each scaffold length, MSE was calculated by comparing to inferences from unfragmented (100 Mb) scaffolds (see methods). Large solid markers and lines show predicted MSE from a linear mixed effect model, with 95% prediction intervals based on simulation.

99 size estimates from all individuals in a given population. Then for each simulated individual, the mean
100 squared error (MSE) was measured between estimates from the fragmented datasets and the loess function
101 for the corresponding population. The MSE was weighted, in discrete time intervals, using the inverse
102 of the sample variance in estimates from the unfragmented datasets (the same individuals as used for
103 the loess fit). This was done to avoid measuring error caused by limited genomic information about the
104 recent and ancient past.

105 Comparisons of the MSE at each fragmentation level (Fig. 1B) suggest that shorter scaffolds do indeed
106 result in less accurate population size estimates. Qualitatively, `mSMC` appears to decline in accuracy at
107 scaffold lengths between 1 Mb and 100 kb for all demographic scenarios, whereas `PSMC` declines in accuracy
108 only in the Recovery scenario, at scaffold lengths between 100 kb and 10 kb.

109 **Mixed effects model**

110 To determine if the observed differences were significant, we fitted a linear mixed-effects model separately
111 for each demographic scenario. The fixed effects were scaffold length and estimation program (`PSMC` vs.
112 `mSMC`), and a random intercept was necessary to account for the repeated measures of each individual
113 at multiple scaffold lengths. Both scaffold length and estimation program were found to be significant
114 predictors of MSE in all demographic scenarios. Two-way interactions between scaffold length terms and
115 estimation program were also significant in all scenarios.

116 **Empirical data**

117 Arguably, the simulated population history scenarios are unrealistic. Simulated data also provides the best
118 possible case in terms of missing data in that there is none. To gauge the impact of using a scaffold-level
119 assembly with real data, we artificially fragmented chromosome 1 from a high coverage human genome,
120 HG00419, a Southern Han Chinese female (The 1000 Genomes Project Consortium, 2015). Population
121 size histories were again estimated using `PSMC` and `mSMC`, for each of the fragmented and unfragmented
122 datasets (Fig. 2).

123 Both programs produced largely the same demographic history when processing long scaffolds, al-
124 though `mSMC` did not estimate population sizes for time intervals as far into the past as `PSMC` (3 Mya vs.
125 10 Mya). For 10 kb scaffold lengths, inferences from `mSMC` are substantially different to those using longer
126 scaffolds, and a small departure is also discernible in the recent past for 100 kb scaffolds. Estimates from
127 `PSMC` have noticeably poorer resolution at the 10 kb scaffold length, but are remarkably consistent for
128 longer scaffolds.

129 The data conversion script provided with `PSMC` (`fq2psmcfa`) ignores scaffolds having fewer than 10000
130 genotype calls by default. This excluded most of the 10 kb scaffolds, due to the presence of one or more
131 missing genotypes. Disabling this filter to retain all scaffolds only marginally improved population size
132 estimates, and only in more ancient time intervals (results not shown). We considered the possibility that
133 with 10 kb scaffolds, `PSMC` might still accurately recapitulate the population size history if provided with
134 more information. To this end chromosome 2 was also partitioned into 10 kb scaffolds and appended

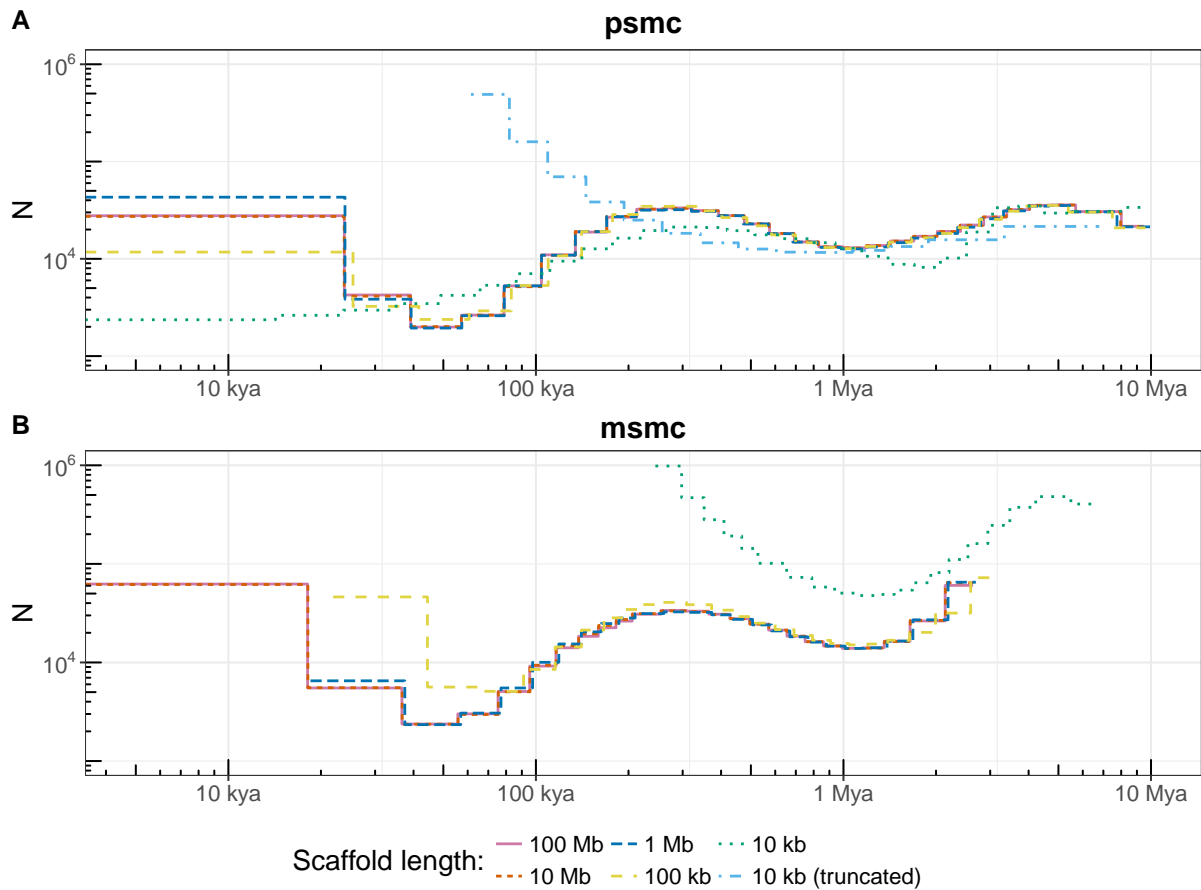


Figure 2: Population size history of HG00419, a Southern Han Chinese individual (The 1000 Genomes Project Consortium, 2015), inferred by **A**) *psmc* and **B**) *msmc*. Empirical data was artificially fragmented to emulate genome sequences aligned to scaffold-level reference assemblies. Population size inferences from *psmc* are consistent down to 100 kb scaffold lengths, with loss of resolution at 10 kb. For *msmc*, stable inferences can be made down to 1 Mb, but accuracy at 100 kb is poor in the recent past, and at 10 kb even broad demographic trends are difficult to discern. Input data to *psmc* for the ‘10 kb (truncated)’ line style had trailing homozygous sites removed from all scaffolds, to match the information content of *msmc* input. Plots were scaled to real time using a 25 year generation time and $1.25e - 8$ mutations per base per generation. kya: thousand years ago; Mya: million years ago;

135 to the chromosome 1 data (doubling the information to ~500 Mb in total). However, the additional
136 information did not alter the result.

137 **msmc discards homozygous tracts at the ends of scaffolds**

138 An input file for `msmc` contains lines that specify the coordinate of a heterozygote site and its distance
139 from the previous heterozygote on the same scaffold. Nothing is specified for coordinates after the last
140 heterozygote, and the scaffold is implicitly truncated here. For short scaffolds this causes substantial
141 information loss. Indeed, short scaffolds may contain no heterozygote sites at all, and input files for such
142 scaffolds are empty.

143 To determine if truncation was a major cause of the different behaviour between `psmc` and `msmc`, we
144 ran `psmc` on 10 kb scaffolds that were artificially truncated to match the information available to `msmc`.
145 Scaffolds containing no heterozygotes were omitted. This output ('10 kb (truncated)' in Fig. 2A), shows
146 a similar trend to that for `msmc` on 10 kb scaffolds, although differences remain.

147 Marginal genealogies with recent coalescent times have accumulated few mutations, so corresponding
148 regions of the genome contain mostly homozygote genotypes. Truncation increases the proportion of
149 heterozygotes, hence recent coalescent times appear older. On short scaffolds, all marginal genealogies
150 are near a scaffold end, so inferences from short truncated scaffolds are more strongly biased to not observe
151 recent coalescent events. Since the population size for each time interval is inversely related to the rate
152 at which pairs of haplotypes coalesce, the smaller number of observations of high homozygosity genomic
153 tracts also means that population size inferences are biased upwards. Both artefacts are noticeable,
154 particularly in the more recent time bins, for `psmc` with artificially truncated 10 kb scaffolds (Fig. 2A)
155 and for `msmc` with 10 kb and 100 kb scaffolds (Fig. 2B).

156 **Conclusion**

157 Reasonable parameter inference in a hidden Markov model relies on observations leading up to, and
158 following, transitions in state. For PSMC, this corresponds to having sufficient sequence contiguity to
159 observe genomic tracts on both sides of historical recombination breakpoints. The chance that a short
160 scaffold will contain a tract covering a recombination breakpoint depends not only on the completeness
161 of the reference assembly, but also the sparsity of breakpoints.

162 Several factors contribute to breakpoint density, including population size, the per base recombination
163 rate, and recombination hotspots. A population suffering a recent and very severe bottleneck will give
164 rise to mostly recent pairwise coalescent times, and few recombination breakpoints, both of which are
165 poorly represented within short scaffolds. Our simulations considered a mammalian recombination rate
166 (3.125×10^{-9} per base per generation) and population size histories that are relevant to many taxa. This
167 suggests that PSMC inference can be reasonable from scaffolds as short as 100 kb for a wide range of
168 datasets.

169 Scaffold level reference assemblies are unlikely to contain equally sized scaffolds, as evaluated here.
170 Generally, a scaffold-level assembly contains tens of long scaffolds and tens of thousands of short scaffolds.

171 In such cases, it is reasonable to exclude scaffolds shorter than 100 kb when running `psmc`, and scaffolds
172 shorter than 1 Mb for use with `msmc`. However, we caution that this guideline may be too optimistic for
173 severely bottlenecked populations or genomic data aligned to a very low quality reference assembly.

174 Materials and Methods

175 Simulations

176 Simulations were performed using `scrm` (Staab *et al.*, 2015), with mutation rate $\mu = 1.25 \times 10^{-8}$ per
177 base per generation and recombination rate $\mu/4$ per base per generation. Simulation output was ar-
178 tificially fragmented during conversion to `psmc` and `msmc` input formats, using a custom Perl script.
179 Demographic inferences were obtained from `psmc` v0.6.5-r67 and `msmc` v1.0.0 for all inputs. Both `psmc`
180 and `msmc` were run with the same time bin parameter (`-p 1*2+15*1+1*2`), although we note that each
181 program calculates time boundaries for the discrete bins differently, so a completely fair comparison is
182 not possible. Scripts used for simulation, format conversion, and running `psmc/msmc` are available from
183 <https://github.com/grahamgower/psmc-error-analysis/>.

185 Mean squared error

186 For each simulated population history scenario and each estimation program, estimates from the unfrag-
187 mented datasets were used to fit a loess function of log population ($\log(N)$) against log time ($\log(t+10)$).
188 The offset of 10 was based on a sensitivity analysis and the smallest non-zero time. An optimal value
189 for the loess smoothing parameter was selected by maximising the corrected AIC (AICc) (Hurvich *et al.*,
190 1998). Mean squared error for individual i in population j was calculated as

$$191 \quad MSE_{ij} = \frac{1}{k} \sum_{m=1}^k (n_{ijm} - \tilde{n}_{.jm})^2 / var_j(m),$$

192 where the sum extends over all k time intervals, n_{ijm} is the log of the population size estimate in interval
193 m , and $\tilde{n}_{.jm}$ is the prediction for the m th time interval from the loess function fitted for the j th population.
194 The variance step function $var_j(m)$ at time interval m , for the j th population, was calculated by splitting
195 time on a log scale into 10 even-width bins and calculating the variance in each bin.

196 Mixed effects modelling

197 Scatter-plots of MSE against scaffold length indicated a cubic relationship between MSE and $\log(\text{scaffold}$
198 $\text{length})$. This was confirmed by comparing residual plots for linear, quadratic, and cubic models. To help
199 numerical consistency of the fitting process, we performed a location scaling of $\log(\text{scaffold length})$.

200 Bivariate analysis of each of the predictors— $\log(\text{scaffold length})$, estimation program, population
201 history scenario, sample ID, and population ID—were used for variable selection. Only $\log(\text{scaffold}$
202 $\text{length})$, estimation program, and population history scenario had a significant relationship with MSE.

203 The linear mixed effects model was fitted using the `lme4` package (Bates *et al.*, 2015) in R (R Core
204 Team, 2017). The fixed effects were $\log(\text{scaffold length})$ and estimation program. Up to two-way inter-
205 action terms were considered for each of the cubic $\log(\text{scaffold length})$ terms with estimation program.
206 To account for repeated measures from each simulated individual due to multiple levels of fragmentation,
207 we included random effects. Both random intercepts and random slopes were considered.

208 All significance testing was performed using the `lmerTest` package (Kuznetsova *et al.*, 2017). All
209 assumptions of the linear mixed-effects models were assessed and regarded as reasonable. The 95%
210 prediction intervals were based on simulation with the `merTools` package (Knowles and Frederick, 2016).

211 Empirical dataset

212 We downloaded the cram alignment file for HG00419, aligned to assembly GRCh38DH, from The 1000
213 Genomes ftp server, and called genotypes with `samtools -q20 -Q20 -C50 ... | bcftools call -c`
214 `...`. The resulting vcf was partitioned into scaffolds of a specific size by modifying the chromosome name
215 and position to which each genotype call corresponded, and was performed separately for each of the
216 scaffold sizes 100 Mb, 10 Mb, 1 Mb, 100 kb, and 10 kb. Input for both `psmc` and `msmc` were filtered to
217 exclude sites with less than half, or greater than double, the mean depth (54.76). The vcf was converted
218 to `psmc` input format with `vcfutils.pl` (distributed with `samtools`) and `fq2psmcfa` (distributed with
219 `psmc`), then `psmc` was run with time bin parameter `-p 4+25*2+4+6`. The same vcf was converted to `msmc`
220 input format with `bamCaller.py` and `generate_multihetsep.py`, both distributed with `msmc-tools`,
221 then `msmc` was run with parameters `-R -p 15*1+15*2`. The time bin parameters for both programs were
222 chosen to be suitable for inferring human demography (Li and Durbin, 2011; Schiffels and Durbin, 2014).

223 Acknowledgements

224 This work was supported by an Australian Government Research Training Program Scholarship [GG],
225 and research fellowships from the University of Adelaide and the Australian Research Council [BL].

226
227 GG, JT, ABR, JS, BL, and NB designed the study. GG performed simulations and processed the
228 empirical data. JT calculated the MSE and performed mixed effects modelling. All authors interpreted
229 the results. GG wrote the manuscript with the help of all coauthors.

230 References

- 231 Bates D, Mächler M, Bolker B, and Walker S. 2015. Fitting linear mixed-effects models using `lme4`. *J*
232 *Stat Softw.* 67(1):1–48.
- 233 Boitard S, Rodríguez W, Jay F, Mona S, and Austerlitz F. 2016. Inferring Population Size History from
234 Large Samples of Genome-Wide Molecular Data - An Approximate Bayesian Computation Approach.
235 *PLoS Genet.* 12(3):e1005877.

- 236 Cleveland W, Grosse E, and Shyu WM. 1992. Local regression models. In: Chambers JM and Hastie
237 TJ, editors, *Statistical Models in S*. Belmont (CA): Wadsworth & Brooks/Cole. p. 309–375.
- 238 Dobrynin P, Liu S, Tamazian G, Xiong Z, Yurchenko AA, Krasheninnikova K, Kliver S, Schmidt-Küntzel
239 A, Koepfli KP, Johnson W, Kuderna LF, García-Pérez R, Manuel Md, Godinez R, Komissarov A,
240 Makunin A, Brukhin V, Qiu W, Zhou L, Li F, Yi J, Driscoll C, Antunes A, Oleksyk TK, Eizirik E,
241 Perelman P, Roelke M, Wildt D, Diekhans M, Marques-Bonet T, Marker L, Bhak J, Wang J, Zhang
242 G, and O'Brien SJ. 2015. Genomic legacy of the African cheetah, *Acinonyx jubatus*. *Genome Biol.*
243 16:277.
- 244 Dutheil JY, Ganapathy G, Hobolth A, Mailund T, Uyenoyama MK, and Schierup MH. 2009. Ancestral
245 population genomics: The coalescent hidden Markov model approach. *Genetics*. 183(1):259–274.
- 246 Feigin CY, Newton AH, Doronina L, Schmitz J, Hipsley CA, Mitchell KJ, Gower G, Llamas B, Soubrier
247 J, Heider TN, Menzies BR, Cooper A, O'Neill RJ, and Pask AJ. 2018. Genome of the Tasmanian tiger
248 provides insights into the evolution and demography of an extinct marsupial carnivore. *Nat Ecol Evol.*
249 2(1):182–192.
- 250 Griffiths RC and Marjoram P. 1997. An ancestral recombination graph. In *Progress in population genetics*
251 *and human evolution*, p. 257–270, New York (NY): Springer.
- 252 Gutenkunst RN, Hernandez RD, Williamson SH, and Bustamante CD. 2009. Inferring the joint de-
253 mographic history of multiple populations from multidimensional SNP frequency data. *PLoS Genet.*
254 5(10):e1000695.
- 255 Harris K and Nielsen R. 2013. Inferring demographic history from a spectrum of shared haplotype lengths.
256 *PLoS Genet.* 9(6):e1003521.
- 257 Hudson RR. 1990. Gene genealogies and the coalescent process. In D. Futuyma and J. Antonovic, editors,
258 *Oxford Surveys in Evolutionary Biology*, volume 7, p. 1–44. New York (NY): Oxford University Press.
- 259 Hurvich CM, Simonoff JS, and Tsai CL. 1998. Smoothing parameter selection in nonparametric regression
260 using an improved Akaike information criterion. *J R Stat Soc Series B Stat Methodol.* 60(2):271–293.
- 261 Knowles JE and Frederick C. 2016. merTools: Tools for Analyzing Mixed Effect Regression Models. R
262 package version 0.3.0.
- 263 Kozma R, Melsted P, Magnússon KP, and Höglund J. 2016. Looking into the past - the reaction of three
264 grouse species to climate change over the last million years using whole genome sequences. *Mol Ecol.*
265 25(2):570–580.
- 266 Kuznetsova A, Brockhoff PB, and Christensen RHB. 2017. lmerTest package: Tests in linear mixed effects
267 models. *J Stat Softw.* 82(13):1–26.
- 268 Li H and Durbin R. 2011. Inference of human population history from individual whole-genome sequences.
269 *Nature.* 475(7357):493–496.

- 270 MacLeod IM, Larkin DM, Lewin HA, Hayes BJ, and Goddard ME. 2013. Inferring demography from
271 runs of homozygosity in whole-genome sequence, with correction for sequence errors. *Mol Biol Evol.*
272 30(9):2209–2223.
- 273 Marjoram P and Wall JD. 2006. Fast “coalescent” simulation. *BMC Genet.* 7:16.
- 274 Mays HL, Hung CM, Shaner PJ, Denvir J, Justice M, Yang SF, Roth TL, Oehler DA, Fan J, Rekulapally
275 S, and Primerano DA. 2018. Genomic analysis of demographic history and ecological niche modeling
276 in the endangered Sumatran rhinoceros *Dicerorhinus sumatrensis*. *Curr Biol.* 28(1):70–76.
- 277 McVean GAT and Cardin NJ. 2005. Approximating the coalescent with recombination. *Philos Trans R*
278 *Soc Lond B Biol Sci.* 360(1459):1387–1393.
- 279 Palamara PF, Lencz T, Darvasi A, and Pe’er I. 2012. Length distributions of identity by descent reveal
280 fine-scale demographic history. *Am J Hum Genet.* 91(5):809–822.
- 281 R Core Team. 2017. R: A Language and Environment for Statistical Computing. R Foundation for
282 Statistical Computing, Vienna, Austria.
- 283 Schiffels S and Durbin R. 2014. Inferring human population size and separation history from multiple
284 genome sequences. *Nat Genet.* 46(8):919–925.
- 285 Sheehan S, Harris K, and Song YS. 2013. Estimating variable effective population sizes from multiple
286 genomes: A sequentially Markov conditional sampling distribution approach. *Genetics.* 194(3):647–662.
- 287 Staab PR, Zhu S, Metzler D, and Lunter G. 2015. scrm: efficiently simulating long sequences using the
288 approximated coalescent with recombination. *Bioinformatics.* 31(10):1680–1682.
- 289 Terhorst J, Kamm JA, and Song YS. 2017. Robust and scalable inference of population history from
290 hundreds of unphased whole genomes. *Nat Genet.* 49(2):303–309.
- 291 The 1000 Genomes Project Consortium. 2015. A global reference for human genetic variation. *Nature.*
292 526(7571):68–74.
- 293 Wiuf C and Hein J. 1999. Recombination as a point process along sequences. *Theor Popul Biol.* 55(3):248–
294 259.
- 295 Zhao S, Zheng P, Dong S, Zhan X, Wu Q, Guo X, Hu Y, He W, Zhang S, Fan W, Zhu L, Li D, Zhang
296 X, Chen Q, Zhang H, Zhang Z, Jin X, Zhang J, Yang H, Wang J, Wang J, and Wei F. 2013. Whole-
297 genome sequencing of giant pandas provides insights into demographic history and local adaptation.
298 *Nat Genet.* 45(1):67–71.